

Roll No. Is : DS5B-2118

Roll No. Is : DS5B-2118

Roll No. Is : DS5B-2118

Roll No. Is : DS5B-2118

Roll No. Is : DS5B-2118

In [1]: !pip install pyspark

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Collecting pyspark
  Downloading pyspark-3.2.1.tar.gz (281.4 MB)
    |████████████████████████████████████████| 281.4 MB 29 kB/s
Collecting py4j==0.10.9.3
  Downloading py4j-0.10.9.3-py2.py3-none-any.whl (198 kB)
    |████████████████████████████████████████| 198 kB 42.5 MB/s
Building wheels for collected packages: pyspark
  Building wheel for pyspark (setup.py) ... done
  Created wheel for pyspark: filename=pyspark-3.2.1-py2.py3-none-any.whl size=281853642 sha256=8b46831582fe33020b51646e24a86fc72a74a84b94240776aece0e2d97207751
  Stored in directory: /root/.cache/pip/wheels/9f/f5/07/7cd8017084dce4e93e84e92efd1e1d5334db05f2e83bcef74f
Successfully built pyspark
Installing collected packages: py4j, pyspark
Successfully installed py4j-0.10.9.3 pyspark-3.2.1
```

```
In [2]: from pyspark.sql import SparkSession
from pyspark.ml import Pipeline
from pyspark.ml.feature import VectorAssembler, StringIndexer, OneHotEncoder
from pyspark.ml.classification import LogisticRegression
from pyspark.ml.evaluation import BinaryClassificationEvaluator

session = SparkSession.builder.appName("HR_Dataset").getOrCreate()
data = session.read.csv("HR comma.csv", header = True, inferSchema = True)
```

In [3]: data.show(10)

```
+-----+-----+-----+-----+-----+
--+-+-----+-----+-----+-----+-----+
|satisfaction_level|last_evaluation|number_project|average_monthly_hours|time_spend_company|Work_accident|left|promotion_last_5years|sales|salary|
+-----+-----+-----+-----+-----+-----+-----+
--+-+-----+-----+-----+-----+-----+
|0.38|0.53|2|157|
3|0|1|0|sales|low|
|0.8|0.86|5|262|
6|0|1|0|sales|medium|
|0.11|0.88|7|272|
4|0|1|0|sales|medium|
|0.72|0.87|5|223|
5|0|1|0|sales|low|
|0.37|0.52|2|159|
3|0|1|0|sales|low|
|0.41|0.5|2|153|
3|0|1|0|sales|low|
|0.1|0.77|6|247|
4|0|1|0|sales|low|
|0.92|0.85|5|259|
```

```

5|          0|    1|          0|sales|    low|          224|
|          0.89|          1.0|          5|
5|          0|    1|          0|sales|    low|
|          0.42|          0.53|          2|          142|
3|          0|    1|          0|sales|    low|
+-----+-----+-----+-----+-----+
--+-+-----+-----+-----+-----+-----+
only showing top 10 rows

```

```
In [4]: data.columns
```

```

Out[4]: ['satisfaction_level',
         'last_evaluation',
         'number_project',
         'average_monthly_hours',
         'time_spend_company',
         'Work_accident',
         'left',
         'promotion_last_5years',
         'sales',
         'salary']

```

```
In [5]: str_idxer = StringIndexer(inputCols = ['sales','salary'], outputCols = ["newsales", "new
```

```
In [6]: one_hot_encoding = OneHotEncoder(inputCols = ["newsales","newsalary"], outputCols = ["ne
```

```
In [7]: vec_ass = VectorAssembler(inputCols = ['satisfaction_level','last_evaluation','number_pr
```

```
In [8]: lr = LogisticRegression(featuresCol= "all_features", labelCol = "left")
```

```
In [9]: mypipeline = Pipeline(stages = [str_idxer, one_hot_encoding, vec_ass, lr])
```

```
In [10]: training, test = data.randomSplit([0.71, 0.29])
```

```
In [11]: lr_model = mypipeline.fit(training)
```

```
In [12]: result = lr_model.transform(test)
```

```
In [13]: result.show(4, truncate = False)
```

```

+-----+-----+-----+-----+-----+
--+-+-----+-----+-----+-----+-----+
-----+-----+-----+-----+-----+
-----+-----+-----+-----+-----+
+
|satisfaction_level|last_evaluation|number_project|average_monthly_hours|time_spend_compa
ny|Work_accident|left|promotion_last_5years|sales          |salary|newsales|newsalary|newsal
es_onehot|newsalary_onehot|all_features                                |rawP
rediction                                |probability                                |prediction
|
+-----+-----+-----+-----+-----+
--+-+-----+-----+-----+-----+-----+
-----+-----+-----+-----+-----+
-----+-----+-----+-----+-----+
+
|0.09          |0.62          |6          |294          |4
|0          |1          |0          |accounting |low          |8.0          |0.0          |(9,
[8],[1.0]) |(2,[0],[1.0]) |(18,[0,1,2,3,4,15,16],[0.09,0.62,6.0,294.0,4.0,1.0,1.0])|
[-0.906113009228541,0.906113009228541] |[0.28779589386338533,0.7122041061366147]|1.0
|
|0.09          |0.62          |6          |294          |4

```

```

|0      |1      |0      |accounting |low      |8.0      |0.0      |(9,
[8],[1.0]) |(2,[0],[1.0]) |(18,[0,1,2,3,4,15,16],[0.09,0.62,6.0,294.0,4.0,1.0,1.0])|
[-0.906113009228541,0.906113009228541] |[0.28779589386338533,0.7122041061366147]|1.0
|
|0.09      |0.77      |5      |275      |4
|0      |1      |0      |product_mng|medium|4.0      |1.0      |(9,
[4],[1.0]) |(2,[1],[1.0]) |(18,[0,1,2,3,4,11,17],[0.09,0.77,5.0,275.0,4.0,1.0,1.0])|
[-0.6706336468375675,0.6706336468375675] |[0.3383549711925885,0.6616450288074115]|1.0
|
|0.09      |0.77      |6      |244      |4
|0      |1      |0      |product_mng|low      |4.0      |0.0      |(9,
[4],[1.0]) |(2,[0],[1.0]) |(18,[0,1,2,3,4,11,16],[0.09,0.77,6.0,244.0,4.0,1.0,1.0])|
[-0.7998272517122431,0.7998272517122431] |[0.31006247261884357,0.6899375273811564]|1.0
|
+-----+-----+-----+-----+-----+
--+-+-----+---+-----+-----+-----+-----+-----+-----+
-----+-----+-----+-----+-----+-----+-----+
-----+-----+-----+-----+-----+
+
only showing top 4 rows

```

```

In [14]: eval = BinaryClassificationEvaluator(rawPredictionCol="rawPrediction", labelCol = "left"

In [15]: eval.evaluate(result)

Out[15]: 0.8297104126919155

```